

1. What is System?

A system is nothing but user of the system and the requirements of users and components that are required to build system

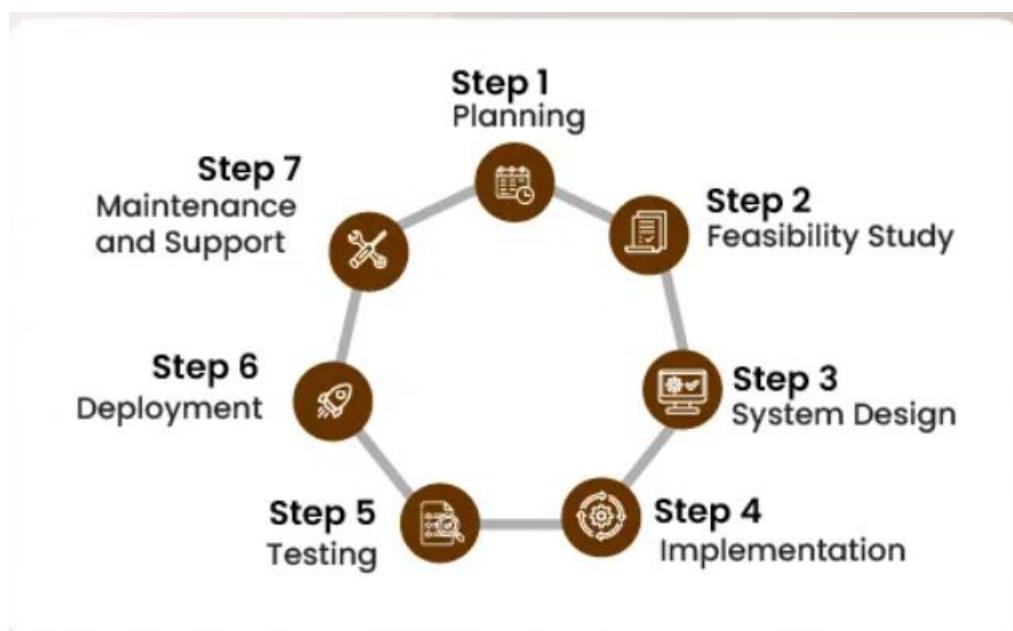
Ex: WhatsApp, Netflix and hot star.

2. What is Design?

Design is a process of understanding the user requirements and selecting the required components like software technologies to serve the need of the system.

3. What is System Design?

Systems Design is the process of defining the architecture, components, modules, interfaces, and data for a system to satisfy specified requirements. It involves translating user requirements into a detailed blueprint that guides the implementation phase. The goal is to create a well-organized and efficient structure that meets the intended purpose while considering factors like [scalability](#), maintainability, and performance.



4. Explain about SDLC?



Stage1: Planning and requirement analysis

Requirement Analysis is the most important and necessary stage in SDLC.

The senior members of the team perform it with inputs from all the stakeholders and domain experts or SMEs in the industry. Planning for the quality assurance requirements and identifications of the risks associated with the projects is also done at this stage. Business analyst and Project organizer set up a meeting with the client to gather all the data like what the customer wants to build, who will be the end user, what is the objective of the product. Before creating a product, a core understanding or knowledge of the product is very necessary.

For Example, A client wants to have an application which concerns money transactions. In this method, the requirement has to be precise like what kind of operations will be done, how it will be done, in which currency it will be done, etc.

Once the required function is done, an analysis is complete with auditing the feasibility of the growth of a product. In case of any ambiguity, a signal is set up for further discussion. Once the requirement is understood, the SRS (Software Requirement Specification) document is created. The developers should thoroughly follow this document and also should be reviewed by the customer for future reference.

Stage2: Defining Requirements

Once the requirement analysis is done, the next stage is to certainly represent and document the software requirements and get them accepted from the project stakeholders.

This is accomplished through "SRS"- Software Requirement Specification document which contains all the product requirements to be constructed and developed during the project life cycle.

Stage3: Designing the Software

The next phase is about to bring down all the knowledge of requirements, analysis, and design of the software project. This phase is the product of the last two, like inputs from the customer and requirement gathering.

Stage4: Developing the project

In this phase of SDLC, the actual development begins, and the programming is built. The implementation of design begins concerning writing code. Developers have to follow the

coding guidelines described by their management and programming tools like compilers, interpreters, debuggers, etc. are used to develop and implement the code.

Stage5: Testing

After the code is generated, it is tested against the requirements to make sure that the products are solving the needs addressed and gathered during the requirements stage. During this stage, unit testing, integration testing, system testing, acceptance testing are done.

Stage6: Deployment

Once the software is certified, and no bugs or errors are stated, then it is deployed. Then based on the assessment, the software may be released as it is or with suggested enhancement in the object segment.

After the software is deployed, then its maintenance begins.

Stage7: Maintenance

Once when the client starts using the developed systems, then the real issues come up and requirements to be solved from time to time.

This procedure where the care is taken for the developed product is known as maintenance.

5. Explain the basic requirements of system design?

System design is a crucial phase in building scalable, efficient, and maintainable systems. The goal of system design is to architect a solution that meets both functional and non-functional requirements while ensuring robustness, scalability, and fault tolerance. Below are the **basic requirements** of system design, covering both **functional** and **non-functional** aspects:

1. Functional Requirements:

These are the core features and behaviours the system must exhibit to fulfil the business goals.

- Core Features:**

- **Data Input & Output:** What type of data will the system handle, and how will it be processed and output? For example, in an e-commerce system, the core feature is the management of products, user accounts, and orders.
- **User Actions & Workflows:** What actions will users take, and what processes should occur behind the scenes? For example, placing an order, checking account balances, or generating reports.
- **External Integrations:** Does the system need to interact with third-party services, databases, or APIs? This could include payment gateways, email services, or cloud storage services.
- **Security & Access Control:** How will users interact with the system (e.g., login, signup)? What types of user roles and permissions will be required? Will there be any encryption for sensitive data?

- Data Management:**

- **Data Flow & Storage:** Where will the data come from (input), where will it be stored (databases), and how will it be accessed? Will you use relational, NoSQL, or file-based storage systems?
 - **Data Processing:** Does the system need batch processing, real-time processing, or a mix? For example, a payment system may need real-time processing, while a reporting system may need batch processing.
-

2. Non-Functional Requirements:

These define the system's qualities or constraints that are not directly related to specific behaviors but are essential to the system's performance, scalability, and usability.

a. Scalability:

- **Vertical Scaling:** Can the system grow by adding more resources (CPU, RAM) to a single server?
- **Horizontal Scaling:** Can the system grow by adding more machines or instances (e.g., microservices, load balancing, distributed databases)?

b. Reliability:

- **Fault Tolerance:** How will the system handle failures? This includes things like **redundancy** (replicating services), **failover mechanisms**, and **graceful degradation** (allowing parts of the system to continue functioning even if other parts fail).
- **High Availability (HA):** Does the system need to be highly available with little downtime? This could mean using geographically distributed servers or cloud services with built-in failover and replication

c. Performance:

- **Low Latency:** How quickly should the system respond to user requests? For example, a search engine should return results in milliseconds.
- **Throughput:** Can the system handle high volumes of transactions or requests? For example, a social media app may need to handle millions of user interactions per second.
- **Response Time:** How fast should the system be able to handle an individual request? Ensure it meets user expectations.

d. Security:

- **Data Encryption:** How will sensitive data (e.g., passwords, financial data) be encrypted both in transit (TLS) and at rest?
- **Authentication & Authorization:** What mechanisms will be used to ensure that users are who they say they are (e.g., OAuth, JWT, SSO)?
- **Data Integrity:** Ensuring that the data cannot be tampered with. This might involve checksums, hash functions, and logging changes.

e. Maintainability:

- **Modularization:** The system should be divided into modules or services that can be easily maintained or updated without affecting the entire system.
- **Code Quality:** Consistent coding standards, proper documentation, and version control (Git) are necessary for long-term maintenance.
- **Testing:** Should include unit tests, integration tests, end-to-end tests, and performance tests to ensure the system works as expected.

f. Usability:

- **User Experience (UX):** How easy and intuitive is it for the user to interact with the system? For example, how intuitive is the design of a web portal or mobile application?
- **Accessibility:** The system should be accessible by people with various disabilities (e.g., screen readers for blind users, color contrast for those with color blindness).

g. Consistency:

- **Consistency Models:** If the system is distributed, what consistency model should be followed (e.g., **strong consistency** vs. **eventual consistency**)? Consider systems like **CAP theorem** (Consistency, Availability, Partition Tolerance).
-

3. Architectural Requirements:

This refers to the overall structure of the system, which will influence how components communicate, are deployed, and scale.

- **Monolithic vs. Microservices:** Will the system be a single monolithic application, or will it be split into microservices? This impacts how the system is scaled, maintained, and deployed.
 - **Client-Server Model:** Will the system follow a traditional client-server model, or will it use modern architectures like serverless or peer-to-peer?
 - **API Design:** How will services communicate? Will there be REST APIs, GraphQL, gRPC, or some other protocol for internal and external communication?
 - **Data Consistency:** If the system is distributed, how will data be synchronized and kept consistent? Options include distributed databases, eventual consistency, or strong consistency mechanisms.
-

4. Data Requirements:

- **Data Redundancy:** Is it important to replicate data for backup purposes, disaster recovery, or to improve availability? Will you need **data partitioning, sharding, or replication**?
 - **Data Modeling:** What data models (relational, document-based, graph) will be most suitable for the system's needs?
 - **Data Governance:** Who owns the data? How is it secured, and what compliance requirements exist (e.g., GDPR, HIPAA)?
-

5. Deployment & Operational Requirements:

These focus on how the system will be deployed, maintained, and monitored in production.

- **Deployment Strategy:** Will the system use cloud-based services (AWS, GCP, Azure), on-premises infrastructure, or a hybrid approach?
 - **CI/CD:** Continuous Integration and Continuous Deployment processes should be in place to enable automated testing, building, and deployment.
 - **Monitoring & Logging:** Tools to monitor system health, performance metrics (Prometheus, Grafana), and centralized logging (ELK stack, Splunk).
 - **Versioning and Rollback:** Can the system's components (especially APIs) be updated safely without downtime? How easy is it to roll back changes?
-

6. Scalability Considerations:

This is particularly important if the system expects to handle a growing number of users, requests, or transactions.

- **Horizontal and Vertical Scaling:** Can the system handle increases in load by adding more instances (horizontal scaling) or by upgrading hardware (vertical scaling)?

- **Load Balancing:** How will traffic be distributed across multiple servers to ensure even distribution and prevent bottlenecks?
 - **Caching:** What strategies will be used to reduce load and improve performance, such as **in-memory caching** (Redis, Memcached) or **content delivery networks (CDN)**?
-

7. Failure Recovery & Disaster Management:

These describe the strategies to handle system downtime, data loss, and other unexpected events.

- **Backup & Recovery:** Ensure data can be restored from backups after a disaster.
 - **Graceful Degradation:** Allow the system to continue working with reduced functionality during partial failures.
 - **Retry Logic and Circuit Breakers:** Implement mechanisms that allow services to retry failed operations and detect cascading failures to prevent system-wide outages.
-

8. Cost Constraints:

- **Cost Efficiency:** Ensure the design adheres to a reasonable budget. Consider cost implications for storage, cloud services, development, and maintenance.
 - **Resource Optimization:** Minimize waste by using auto-scaling, serverless computing, or efficient data storage solutions.
-

Conclusion:

System design requires a thorough understanding of both the **functional** and **non-functional** requirements to create a well-balanced and effective system. While functional requirements focus on the core business logic, non-functional requirements ensure the system can scale, perform well, remain secure, and be maintainable over time. A good system design involves a balance between these two areas to ensure that the system not only meets business goals but also performs reliably and efficiently in production.

6. Explain the components of system design?

Components can be divided into two parts

1. Logical Entities
2. Tangible Entities(Technology)

Logical Entities: Logical entities are made up of databases, applications, communication protocols (tcp/udp etc), https requests like REST, Soap. All these components run on a system or platform and actually come under system requirements. To be precise things like Cache, message Queues, infrastructure etc.

Tangible Entities: Tangible entities are those which provide bulk to the logical entities. They are software/databases etc which fulfills the above. Eg:- Redis, mysql, RabbitMq, API's etc.

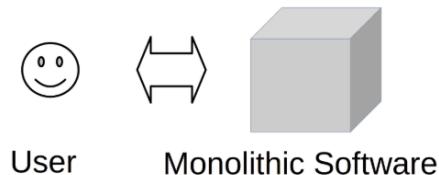


3. Explain different system design architectures?

- ✓ **Layered Architecture:** divides the system into logical layers (display, business logic, and data), each with its responsibilities. It improves modularity and eases maintenance but may result in dependencies between levels.

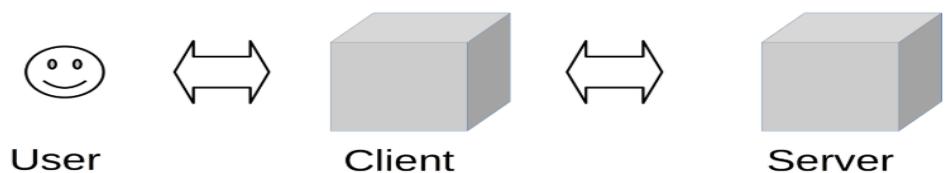
One-Tier (Monolithic) Architecture

Referred to as monolithic architecture, the One-Tier system consolidates all components, spanning the user interface to data storage, within a single executable or process. Predominantly employed for simpler applications, this architecture streamlines development and deployment processes:



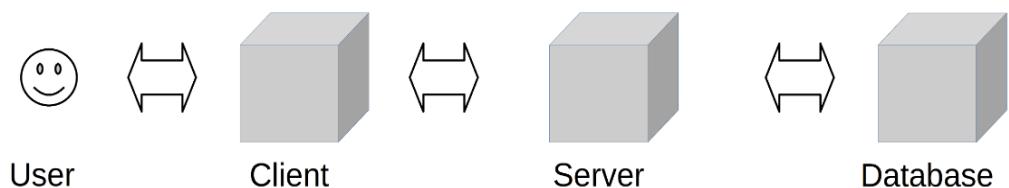
Two-Tier (Client-Server) Architecture

The Two-Tier architecture bifurcates the application into two fundamental components: the client and the server. The client assumes responsibility for managing the user interface, while the server handles crucial functions such as data storage and business logic. Ideal for medium-sized applications, this architecture optimizes performance and resource allocation:



Three-Tier (Data-Server) Architecture

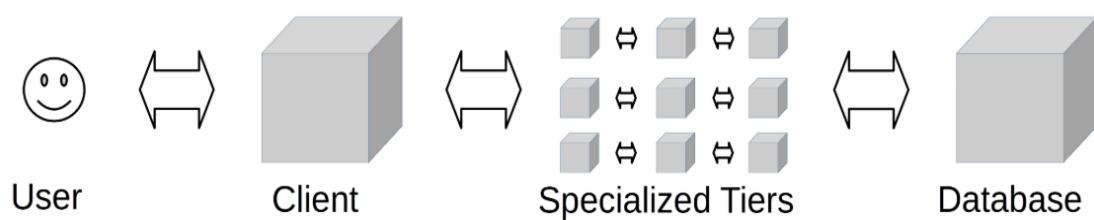
Extending the structural hierarchy, the Three-Tier architecture segments the application into presentation, application, and data tiers. This segregation enhances scalability and maintainability by effectively separating concerns. Well-suited for complex systems, this architecture ensures a more organized and adaptable development environment:



N-Tier (Information-Server) Architecture

N tier architecture also referred to as distributor architecture, where n stands for the number of tiers. The difference between the 3 tier and n tier is that there is more than one application server intermediatelying between the user interface layer and the database layer.

This is done to distribute the business logics among different servers so that all of them can perform individually. This architecture is very flexible with the increase in the number of databases server. The n tier architecture based application is more secure than the other tiers of architecture.

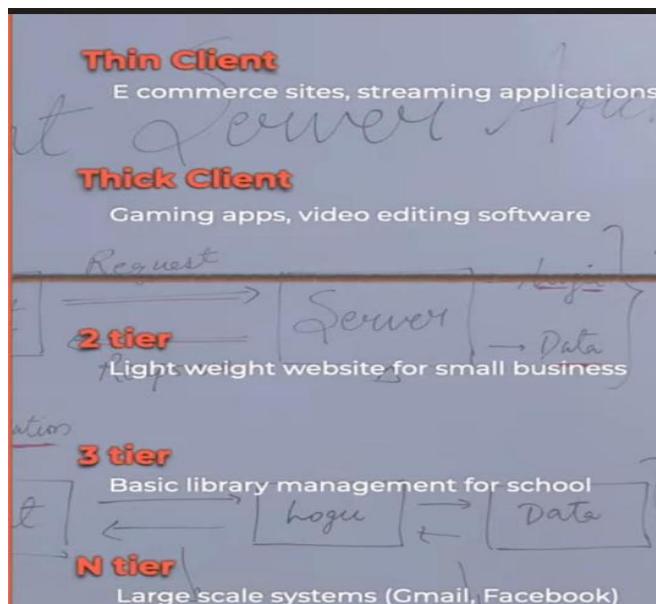


- ✓ **Client-Server Architecture:** This architecture divides the system into client and server components, with clients requesting services or resources from servers. It allows for scalability and centralized control, but it necessitates network connectivity.
- ✓ **Model-View-Controller (MVC):** The program is divided into three interconnected parts by this architecture: the Model (information and rationale), the View (show), and the Regulator (which oversees client input). It encourages modularity, testability, and reuse of code.
- ✓ **Microservices Architecture:** Break down the application into smaller, self-contained services. Each service focuses on a unique business feature, allowing for greater scalability, flexibility, and ease of maintenance but increasing the complexity of managing dispersed systems.
- ✓ **Event-driven architecture (EDA)** focuses on creating, detecting, consuming, and responding to events. Components interact via events, allowing loose coupling, scalability, and responsiveness.
- ✓ **Service-Oriented Architecture (SOA):** divides the system into well-defined, self-contained, and loosely connected services. It encourages reusability, flexibility, and interoperability but can add complexity to dependency management.

4. Difference between thick and thin clients?

Thick Client: If logic and processing sits and client side is known as Thick Client. Ex: Outlook

Thin Client: If logic and processing sits and server side is known as Thick Client. Ex: Netflix



5. Explain about proxy and its types and disadvantages?

A proxy is essentially an intermediary that sits between a client and a server, forwarding requests and responses between the two. It can be used for various purposes such as load balancing, security, caching, and more.

Real-World Example:

Imagine a bookstore where some books are stored in the main shop, and others are kept in a warehouse. As a customer, you don't have direct access to the warehouse. Instead, when you request a book that isn't on the shelves, a staff member (the intermediary) goes to the warehouse, retrieves the book, and hands it to you. This staff member acts as a proxy between you (the client) and the warehouse (the server), ensuring that the system fulfills the request efficiently without you needing to navigate through the entire storage facility.

Types of Proxies:

Forward Proxy: A forward proxy sits between the client and the server, forwarding the client's requests to the server and returning the server's responses to the client. It essentially acts on behalf of the client, masking the client's identity from the server.

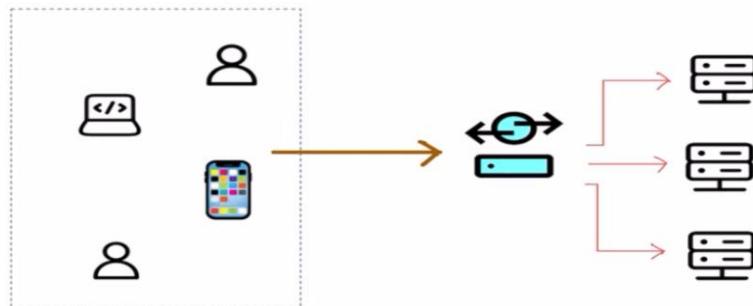


Fig: A forward proxy acts on behalf of clients

Benefits of Forward Proxy:

Privacy & Anonymity: The forward proxy hides the client's identity. This is useful for privacy, as the server won't know who exactly is making the request.

Example: When you browse the internet using a proxy, websites won't see your IP address. Instead, they will see the proxy's IP.

Content Filtering & Caching: The proxy can filter requests, blocking access to certain content or caching frequently requested items, reducing load times.

Example: In a bookstore, the staff member (proxy) could know that certain books are in high demand and retrieve them in advance, ensuring faster service.

Security: Forward proxies can add an extra layer of security, protecting clients from malicious servers by filtering harmful content.

Reverse Proxy: A reverse proxy sits in front of the server, intercepting client requests and forwarding them to the appropriate server. It acts on behalf of the server, masking the server's identity from the client.

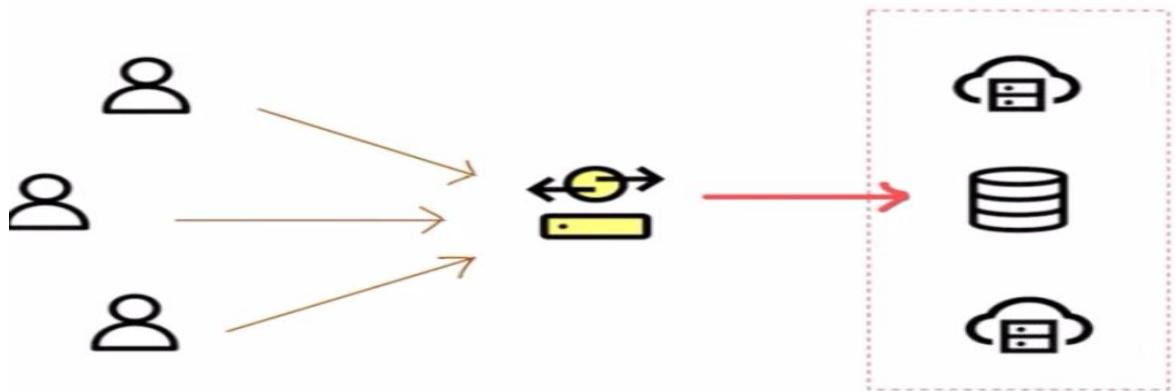


Fig: Reverse Proxy works on behalf of the servers

Benefits of Reverse Proxy:

Load Balancing: It can distribute incoming requests across multiple servers to prevent any single server from being overloaded.

Security: Reverse proxies mask the identity of the server, adding a layer of protection. Clients never know the true IP address of the server, reducing the risk of direct attacks.

SSL Termination: Reverse proxies can handle encryption and decryption for SSL (HTTPS) traffic, reducing the load on the backend servers.

Working from Home Example: When you work from home, you often connect to your company's network through a reverse proxy. This routes your requests to the correct server within the company's internal network, ensuring that your communication is secure and that external users can't directly access the company's servers.

Challenges of Using Proxies

While proxies provide numerous benefits, they also introduce some challenges.

- **Latency:** Adding a proxy introduces an additional step in the communication process. This can lead to increased latency if not optimized properly.
- **Single Point of Failure:** If the proxy goes down, it can bring the entire system to a halt, as both clients and servers depend on it to relay messages.
- **Complex Configuration:** Proxies can add complexity to the overall system design, requiring careful configuration and maintenance.

6. Difference between proxy and gateway?

A **Proxy** and a **Gateway** are both network intermediaries that mediate traffic between clients and servers, but they serve different purposes and have distinct use cases. Below is a detailed breakdown of the differences between a **Proxy** and a **Gateway**:

- **Proxy:**

- A **proxy** acts as an intermediary between a client and a server. It typically forwards requests from clients to the appropriate servers and can return the server's response to the client.
 - Proxies are often used for purposes like security (hiding client identity), load balancing, and content caching.
 - **Gateway:**
 - A **gateway** is a more sophisticated intermediary that not only forwards requests but also manages and translates communication between different protocols or services.
 - Gateways often deal with converting data formats, protocol translation, or API routing.
-

Functionality:

- **Proxy:**
 - **Request forwarding:** A proxy forwards client requests to the server.
 - **Caching:** Can cache responses to improve performance and reduce server load.
 - **Anonymity:** Can hide the client's IP address from the server.
 - **Load Balancing:** Distributes incoming requests across multiple servers to balance the load.
 - **Security:** May block certain types of traffic or filter content.
- **Gateway:**
 - **Protocol Translation:** Can translate between different protocols (e.g., HTTP to MQTT, SOAP to REST).
 - **Routing:** Routes requests to different microservices, potentially transforming the request format.
 - **API Aggregation:** A gateway can aggregate data from multiple services into a single response.
 - **Security:** Provides authentication, authorization, and other security policies at a higher level, typically at the service-level (e.g., API Gateway).
 - **Request and Response Transformation:** Modifies requests and responses to meet the required formats or protocols for the downstream services.

7. Difference between Vertical and Horizontal scaling?

| Feature | Vertical Scaling | Horizontal Scaling |
|-----------------|--|---|
| Definition | Increasing the power of a single machine by adding resources (CPU, RAM, etc.). | Increasing capacity by adding more machines or nodes. |
| Scalability | Limited by the hardware capabilities of a single machine. | Almost unlimited, as new machines can be added as needed. |
| Cost | Initially cheaper but becomes more expensive as the hardware upgrade limits approach. | Can be more cost-effective at large scale due to distributed load, but requires more infrastructure. |
| Simplicity | Easier to set up and manage as it involves fewer components (only one server). | More complex to manage, as it involves handling multiple servers and ensuring synchronization. |
| Fault tolerance | Single point of failure — if the machine goes down, the entire system can go down. | High fault tolerance — if one machine fails, others continue running. |
| Performance | Can hit a limit on performance due to hardware constraints. | Better suited for high-demand systems that require scaling beyond the capabilities of a single server. |
| Use cases | Suitable for applications that don't require distribution or where scaling a single machine is sufficient. | Ideal for web applications, cloud systems, and distributed services that need to handle large amounts of data or traffic. |
| Maintenance | Easier to maintain since you only have one server to manage. | More complex to maintain because of distributed nodes and potential data consistency challenges. |

8. Explain the best place to implement the caching in micro services?

The **best place to implement caching** in a microservices architecture depends on the type of data being cached, performance goals, and architectural considerations. Caching can be implemented at multiple layers, each serving a specific purpose. Here's a breakdown of optimal caching layers in a microservices architecture:

1. API Gateway Level (Edge Caching)

Caching implemented at the API Gateway level to store frequently requested responses before requests hit the microservices. It is best for Static content, public API responses, authentication tokens

2. Service Level (In-Memory Cache)

Each microservice caches data internally to reduce database queries. Caching frequently queried data such as user sessions, authentication details, or product catalogues.

3. Distributed Cache (Between Services)

A shared cache across multiple instances of a microservice. Caching commonly shared data like product information or configuration settings across multiple microservices.

4. Database Level (Query Result Caching)

Caching query results or entire datasets directly at the database level. Caching frequently executed read-heavy queries (e.g., product details, user profiles).

5. Content Delivery Network (CDN) Caching for Static Content

Caching static assets like images, CSS, JavaScript files at globally distributed servers. Caching static content for faster global delivery.

6. Client-Side Caching (Browser Cache)

Caching data at the client-side for reducing redundant network requests. Caching UI state, static assets, and API responses.

7. Cache for Event-Driven Systems (Message Caching)

Caching messages in a distributed messaging system to avoid reprocessing. Event sourcing, stream processing.

9. Difference between Orchestration & Choreography design patterns?

Orchestration

In the realm of microservices architecture, orchestration refers to a centralized management pattern where one service, the orchestrator, directs the interactions between microservices. Picture an orchestra conductor: he doesn't play each instrument, yet he controls the entire performance by guiding every musician when to start, stop, and adjust their pace. The same applies to the orchestration approach: the orchestrator service doesn't perform the tasks but dictates what individual services should do and when. This approach grants the team a high level of control and oversight, particularly useful in complex workflows.

Microservices Orchestration Example

Consider an e-commerce application as an example. When a customer places an order, multiple tasks need to happen sequentially: validate the product availability, charge the customer's credit card, update the inventory, and finally confirm the order to the customer. To manage these tasks, an orchestrator service takes the lead.

1. It first requests the inventory service to check product availability.
2. Upon confirmation, it then communicates with the payment service to process the customer's credit card payment.
3. If the payment goes through, it updates the inventory.
4. Lastly, it alerts the customer service to send an order confirmation to the customer.

All these actions follow specific steps and sequences under the watchful eye of the one orchestrator. This is a classic example of the orchestration pattern, where a central component manages and coordinates the actions of various services. This approach ensures the tasks are executed in the correct sequence, eliminating the risk of processing a payment for an unavailable product.

What is Choreography

In contrast to orchestration, choreography designates a decentralized pattern in microservices architecture. If we stick to the arts analogy, you can think of it like a dance

troupe. There is no one person dictating every move, but each dancer knows what steps to take based on the music and other dancers' actions. Similarly, in choreography design, no central process directly controls the actions of each service. Instead, each microservice is autonomous and decides its actions based on the events or the actions of other services. This design encourages individual services' independence and can be advantageous where there is a need for high flexibility and adaptability.

Microservices Choreography Example

Imagine a ride-sharing application like Uber or Lyft. When a ride is booked, many independent microservices need to work together. Here's how choreography can be applied:

- Order Service emits an Order Placed event.
- Payment Service listens and processes payment, then emits a Payment Completed event.
- Inventory Service listens for the payment event and updates stock, emitting an Inventory Updated event.
- Shipping Service listens and initiates shipping once the inventory is updated.

In this scenario, there's no central controller dictating the flow of operations. Instead, each service acts on its own and interacts with adjacent ones as required. This is a typical example of choreography's decentralized control, where each microservice performs its part in harmony to ensure a smooth ride.

10. Explain about Service Discovery?

Service Discovery in microservices architecture is a mechanism that allows services to dynamically discover and communicate with each other without hardcoding the service locations. It helps in managing the complexity of distributed systems where services can scale dynamically.

Client-Side Discovery:

- The client is responsible for discovering the service and selecting an instance to send the request.
- **Flow:** Client → Service Registry → Service Instance
- **Tools:** Netflix Eureka, Consul
- **Example:** Load balancers in the client fetching available service instances and routing traffic directly.

Server-Side Discovery:

- The client makes a request to a **load balancer** or **API gateway** which fetches service instances from the service registry and forwards the request.
- **Flow:** Client → Load Balancer/API Gateway → Service Registry → Service Instance
- **Tools:** AWS ELB, Kubernetes Ingress, Traefik

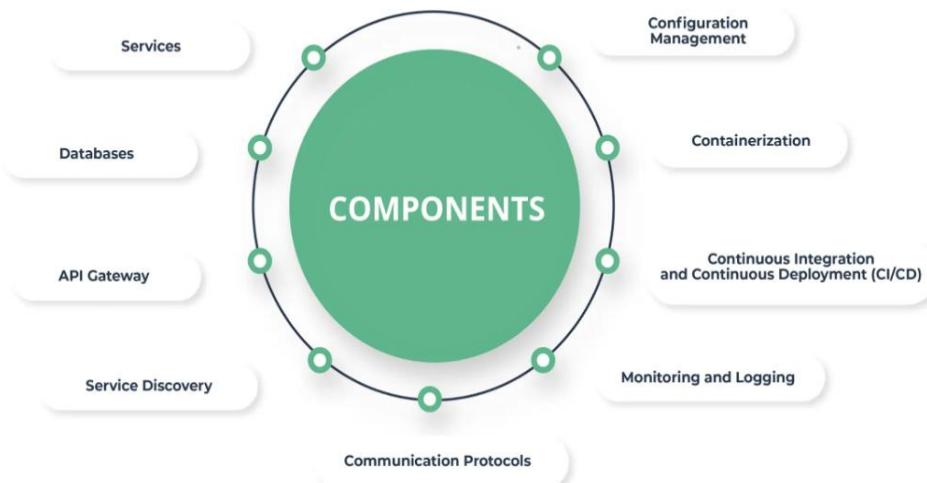
11. Explain about Load Balancing?

Load balancing is the process of disseminating network traffic on the number of servers in the network. This makes certain that no single server suffers from over load demand which enhances reliability and performance. Load balancers help in managing traffic spikes, preventing server overloads, and ensuring high availability by rerouting traffic from failed servers to healthy ones.

There are several types of load balancers:

- **Hardware Load Balancers:** Hardware appliances designed for traffic distribution, which is necessary in large-scale organizations where high speed and stability of the network are essential. These include F5 and Citrix NetScaler.
- **Software Load Balancers:** Load balancers like Nginx or HAProxy or any other application or service that routes traffic like Envoy. These can be installed on normal web servers and are flexible and scalable.
- **Cloud-based Load Balancers:** Cloud services in the form of managed services from cloud providers such as Amazon Web Service Elastic Load Balancer or Google Cloud Load Balancer. These are relatively simple to configure and the load can be easily balanced and distributed as per the traffic.

12. Explain about micro services components?



1. Services

The core component of microservice architecture is the services themselves. Each service is designed to accomplish a specific task or function and can operate independently of the others.

Services are loosely coupled, meaning changes to one service should not directly impact the functionality of another. They encapsulate specific business logic and data and can be developed, deployed, and scaled independently.

2. Databases

Each service typically manages its database in a microservice architecture, allowing for data autonomy and ensuring that database dependencies do not hinder its operations.

This approach, known as database per service, helps isolate the services, improve fault tolerance, and allow the use of different database technologies suited to each service's needs.

3. API Gateway

The API Gateway is the entry point for clients to access the application's services. It is responsible for routing requests to the appropriate microservice, aggregating the results, and returning them to the requester.

The API Gateway can also handle cross-cutting concerns such as authentication, SSL termination, and rate limiting, thereby simplifying the client's interaction with the microservices.

4. Service Discovery

Service discovery is a critical component in microservices architecture, enabling services to discover and communicate with each other dynamically.

This is essential because, in a [cloud](#) environment, services can frequently change locations and IP addresses. Service discovery mechanisms ensure that these services can locate and communicate with one another reliably.

5. Communication Protocols

Communication between microservices is a fundamental aspect of microservice architecture. This communication can be synchronous (e.g., HTTP, REST, gRPC) or asynchronous (e.g., messaging queues like RabbitMQ or Apache Kafka).

The choice of communication protocol depends on the application's specific requirements, such as response time, reliability, and scalability.

6. Configuration Management

Configuration management tools help manage the configurations of the various centralized services. This is crucial for maintaining consistency across environments and simplifying the management of service configurations, especially as the number of services grows.

7. Containerization

Containerization technologies like Docker and Kubernetes play a significant role in microservices architecture.

They provide a lightweight, consistent environment for deploying and running services, making it easier to manage lifecycle issues, scale services, and ensure that the application runs the same way in development, testing, and production environments.

8. Continuous Integration and Continuous Deployment (CI/CD)

CI/CD practices are integral to microservices architecture, enabling frequent and reliable delivery of individual service updates.

Continuous integration involves automatically testing and building services upon code changes, while continuous deployment automates the deployment of services to production. This facilitates a more agile development process and quicker release cycles.

9. Monitoring and Logging

Given the distributed nature of microservices, monitoring, and logging are essential for observing the health of services, diagnosing issues, and understanding the application's behavior.

Tools and platforms that offer centralized logging and application performance monitoring (APM) are commonly used to manage these tasks.

13. Explain about Gateway role in your project?

In our project will APIGEE gateway for authorization and load balancer to serve the consumer and producer needs. These the steps will follow in our project,

- Producer will onboard his application and its load balancing application in APIGEE gateway portal.
- Consumer will onboard subscriber application with the scopes provide by the provider. Provider has to approve the consumer subscription request.
- Producer will provide the application key specific to each consumer to access producer services.
- As a initial step consumer will get the access token from producer with the help of application key provided by the producer.
- If application key is valid one then consumer will provide the access token otherwise 403 forbidden error will be returned. Generated access token has expiry time like 4 hours.
- At the time of producer application onboarding producer will add his actual application load balancing URL at APIGEE gateway.
- Consumer will invoke the producer with the newly generated access token via APIGEE gateway load balancing URL. Token authorization will happened at APIGEE gateway itself.

14. How to build fault tolerant/high available micro service application?

Creating highly available microservices requires a combination of good design, deployment, and monitoring practices. Here are some best practices to consider:

Service Replication: Deploy multiple instances of each microservice across availability zones or regions.

Independent Services: Each microservice should operate independently to reduce dependencies. This ensures that issues in one service do not affect others.

Stateless Services: Design services to be stateless so they can be easily scaled up or down and replaced without losing any information.

Use of Containers: Deploy microservices in containers (e.g., Docker) to ensure consistency across different environments and to simplify deployment and scaling.

Service Registry: Implement a service registry to keep track of available services and their instances.

Load Balancing: Use load balancers to distribute traffic evenly across multiple instances of your services.

Automated Scaling: Set up automated scaling based on demand to handle traffic spikes and ensure consistent performance.

Health Checks: Regularly perform health checks to detect and recover from failures quickly.

Centralized Logging and Monitoring: Implement centralized logging and monitoring to keep track of service performance and quickly identify and troubleshoot issues.

Circuit Breaker Pattern: Use the circuit breaker pattern to prevent cascading failures by stopping operation calls to failing services.

Distributed Tracing: Implement distributed tracing to track requests across different services and understand system behaviour.

Database Replication: Use database replication strategies to ensure data availability

Master-Slave Replication: For read-heavy workloads.

Multi-Master Replication: For write-heavy workloads.

Circuit Breaker Pattern: Prevent cascading failures by temporarily stopping requests to a failing service. Use tools like Netflix Hystrix or Resilience4j.

Retry with Backoff: Implement retry mechanisms with exponential backoff to handle transient failures without overloading the system.

Event-Driven Architecture: Decouple services using asynchronous messaging or event-driven patterns (e.g., with Kafka, RabbitMQ).

15. Explain the best practices to scale the micro services?

Scaling microservices effectively ensures that your application can handle increased loads while maintaining performance and reliability. Here are some best practices:

Horizontal Scaling: Add more instances of your microservices to distribute the load rather than scaling a single instance vertically.

Auto-scaling: Set up auto-scaling to automatically add or remove instances based on defined metrics such as CPU usage or request rates.

Load Balancing: Use load balancers to evenly distribute traffic across multiple service instances.

Asynchronous Processing: Use message queues and background processing to handle workloads asynchronously, reducing the load on your microservices.

Lazy Load and Cache: Implement caching and lazy loading to reduce the load on backend services by storing frequently accessed data.

Database Sharding: Distribute database load by sharding data across multiple databases.

Global Distribution: Deploy services across multiple regions to reduce latency and improve performance for users across different geographies.

16. How to handle the rate limits in micro services?

Handling rate limits in highly used applications ensures that your services remain performant and available even under heavy load. Here are some strategies to effectively manage rate limiting:

Implement Proper Rate Limits: Set appropriate rate limits based on your application's traffic patterns and user needs. You can apply different rate limits for various endpoints or based on user roles.

Use an API Gateway: Employ an API gateway to manage rate limiting across multiple services at a central point. API gateways like Kong, Apigee, or AWS API Gateway provide built-in rate limiting features.

Caching: Cache responses for frequently accessed endpoints to reduce the number of redundant requests your server needs to handle.

Graceful Degradation: Design your application to degrade gracefully by providing limited functionality or cached responses when rate limits are hit.

Communicate Rate Limits: Document and communicate your rate limits to users in your API documentation to set clear expectations.

API Keys and User Authentication: Use API keys and user authentication to apply rate limits based on individual users or clients.

Fallback Mechanisms: Implement fallback mechanisms, such as returning cached responses or default values, when rate limits are exceeded.

Visualization: Use tools like Grafana or Kibana to visualize rate-limiting patterns.

17. Difference between load balancer and Gateway?

Gateway:

- Managing APIs and microservices in a distributed architecture.
- Implementing a single entry point for mobile, web, and IoT clients.
- Adding cross-cutting concerns like authentication, logging, or rate limiting.

Load Balancer:

- Scaling web servers or application servers to handle increased traffic.
- Distributing traffic among multiple instances of a single service.
- Improving fault tolerance and high availability by rerouting traffic from failed servers.

18. Can we use load balancer and gateway together and its benefits?

Absolutely! You can, and often should, use load balancers in conjunction with gateways. Combining them allows you to take advantage of each component's strengths, resulting in a more robust and efficient architecture.

Common Setup

API Gateway at the Front: The API gateway sits at the edge of your system, acting as the entry point for all client requests.

Load Balancer Behind the Gateway: Load balancers are positioned behind the gateway to distribute traffic across multiple instances of your services.

Workflow:

Client Requests: Users send their requests to your application.

Gateway Handling: The API gateway first processes these requests. It handles authentication, authorization, rate limiting, and request transformation.

Traffic Distribution: Once the gateway processes a request, it forwards the request to the load balancer.

Load Balancer Work: The load balancer then distributes the request to one of the available service instances, ensuring even traffic distribution and handling failover.

Benefits

Enhanced Security: The gateway handles security concerns such as authentication and authorization before traffic reaches the internal network.

Efficient Load Management: Load balancers ensure that no single service instance is overwhelmed, improving availability and reliability.

Scalability: Combining both allows you to scale your services horizontally and handle increased load efficiently.

Centralized Monitoring and Logging: The gateway can provide centralized logging and monitoring, giving visibility into request flows and potential bottlenecks.

19. Explain about database scaling?

Scaling databases is a crucial aspect of ensuring your application can handle increased load and maintain performance. There are two main strategies for database scaling: **vertical scaling** and **horizontal scaling**.

Vertical Scaling (Scaling Up): Vertical scaling involves increasing the power of a single database server by adding more CPU, RAM, or storage. This approach is straightforward but has limitations, as there is a maximum capacity for a single machine.

Horizontal Scaling (Scaling Out): Horizontal scaling involves adding more servers to distribute the load across multiple machines. This can be achieved through techniques like sharding and replication.

Sharding: Sharding divides your dataset into smaller pieces, called shards, and distributes them across multiple servers. Each shard is responsible for a subset of the data. Sharding allows you to scale out by adding more servers as your data grows.

Replication: Replication involves copying data across multiple database servers. There are different replication strategies, such as master-slave (primary-secondary) and multi-master.

Master-Slave Replication

Master/Primary: Handles all write operations.

Slaves/Secondaries: Serve read operations and replicate data from the master.

Multi-Master Replication: In multi-master replication, multiple servers handle both reads and writes, replicating data among themselves.

20. Explain about sharding?

Sharding is a database scaling technique where a large dataset is partitioned into smaller, more manageable pieces called shards. Each shard is hosted on a separate database server, allowing for horizontal scaling and improved performance. Sharding enables databases to handle large volumes of data and high-transaction workloads by distributing the load across multiple servers.

Sharding Strategies:

Range Sharding: Data is divided into ranges based on the sharding key. Each shard contains a specific range of data. This is useful for scenarios where data is naturally ordered, such as time-series data.

Hash Sharding: A hash function is applied to the sharding key to determine the shard where the data should reside. Hash sharding ensures an even distribution of data, preventing hotspots where some shards have significantly more data than others.

Geographic Sharding: Data is divided based on geographic regions. Each shard contains data relevant to a specific region. This approach is commonly used in applications with geographically distributed users.

Directory-Based Sharding: A lookup table (directory) is used to map the sharding key to the corresponding shard. This allows for more flexibility in determining shard placement but adds complexity in maintaining the directory.

Example Use Cases:

E-commerce Platforms: Handling large volumes of transactions and product data by distributing the load across multiple shards.

Social Media Applications: Managing user data and interactions for millions of users by sharding based on user ID or geographic region.

Financial Systems: Ensuring high performance and availability for transaction processing and account management by sharding based on account number or region.

Example: MongoDB, Cassandra etc.

21. How to handle the database failures in distributed architecture?

Data Replication: Replicate your data across multiple nodes or data centres. Use techniques like master-slave (primary-secondary) replication or multi-master replication to ensure that copies of your data are always available.

Failover Mechanisms: Implement automatic failover processes to switch to a replica node if the primary node fails. Database systems like MySQL, PostgreSQL, and MongoDB have built-in support for automated failover.

Regular Backups: Perform regular backups of your databases to ensure that you have recent data to restore in case of a major failure. Use automated backup solutions to streamline this process.

Health Checks: Implement continuous health monitoring to detect failures early. Use automated health checks to monitor database nodes and alert your team to potential issues.

Centralized Logging: Use centralized logging systems to track database operations, detect anomalies, and troubleshoot issues quickly.

Circuit Breaker Pattern: Use the circuit breaker pattern to stop requests to a failing database and prevent cascading failures. This pattern helps to isolate and limit the impact of failures.

Graceful Degradation: Design your application to degrade gracefully by providing limited functionality or cached responses when the database is unavailable. This helps maintain a basic level of service during failures.

Sharding: Partition your database into smaller shards to distribute the load and minimize the impact of a failure on a single shard. This also improves performance and scalability.

Load Balancers: Use load balancers to distribute traffic evenly across database nodes and reduce the risk of overloading a single node.

Caching: Implement caching mechanisms to reduce the load on your database by serving frequently accessed data from the cache.

22. Explain cache evict process in distributed architecture?

Cache eviction is the process of removing data from the cache to make room for new data when the cache reaches its storage limit. In a distributed architecture, effective cache eviction policies are vital for maintaining performance and ensuring that the most relevant data is kept in the cache. Here are some common cache eviction policies and strategies used in distributed systems:

Least Recently Used (LRU): Evicts the least recently accessed items first. This policy assumes that data accessed recently is likely to be accessed again soon.

Least Frequently Used (LFU): Evicts the least frequently accessed items first. This policy is based on the assumption that items accessed less frequently are less valuable.

First In, First Out (FIFO): Evicts the oldest items first, regardless of how often they are accessed. This policy can be useful when data has a natural expiration time.

Random Replacement (RR): Evicts items randomly. This method can prevent certain pathological access patterns from causing unfair eviction of other items.

Time-to-Live (TTL): Sets an expiration time for each item, evicting it after its TTL has elapsed. This is useful for ensuring that stale data is automatically removed from the cache.

23. Explain best protocols for micro services?

gRPC: Best for high-performance, low-latency communication (e.g., backend-to-backend services, real-time applications).

HTTP/HTTPS (REST): Best for ease of integration, wide adoption, and public APIs (e.g., web services, third-party integrations).

Kafka: Best for event-driven architectures, real-time data streaming, and decoupling services (e.g., event sourcing, analytics).

RabbitMQ: Best for reliable messaging and queuing (e.g., background tasks, job queues).

WebSockets: Best for real-time communication with persistent connections (e.g., live updates, chat services).

MQTT: Best for IoT devices, telemetry, and constrained environments.

24. In distributed architecture who will select the leader in case of primary node failure in MongoDB?

In MongoDB, leader election is managed by a process called **replica set election**. When a primary node fails, MongoDB automatically selects a new primary from the available secondary nodes.

Detection of Primary Failure: Secondary nodes detect the failure of the primary through heartbeat checks. Each secondary node sends a heartbeat to the primary at regular intervals. If a secondary node stops receiving heartbeats from the primary node for a certain period (typically 10 seconds), it concludes that the primary has failed.

Election Initiation: Once a secondary node detects that the primary is down, it starts the election process. MongoDB uses a replica set election protocol (based on Paxos consensus) to elect a new primary. Eligible candidates for the primary role are the secondaries, and sometimes arbiters can be used for tie-breaking.

Voting: Each node in the replica set (except for the failed primary) can vote during the election. A node is eligible to vote if it has an up-to-date copy of the data and is not currently in a read-only state (i.e., it is an eligible secondary). The voting process is based on the "most recent operation time" (Oplog), which means that the secondary with the most recent replication state may be preferred.

Election Criteria: A new primary is selected when a majority of the nodes (more than half of the replica set members) vote for the same secondary node. The node with the highest Oplog entry (the most recent data update) is usually preferred because it is considered the most up-to-date.

New Primary Election: The node that receives the majority of the votes becomes the new primary node. The other secondaries in the replica set will start replicating from the newly elected primary.

Replication: The new primary starts accepting write operations, and the secondaries replicate from it. The system is now back to normal operation with a new primary, and the former primary (if it comes back online) may be automatically re-elected as a secondary node (or a new secondary may be chosen to replace it, depending on the state).

25. In distributed architecture who will select the leader in case of primary node failure in Kafka?

In Apache Kafka, the process of selecting a **new leader** in case of a **primary node failure** (more precisely, a **partition leader failure**) is managed by **Kafka's Zookeeper** or, in newer versions, by **KRaft mode** (Kafka Raft Protocol). The leader election is specific to **Kafka topics** and **partitions**, and it's a key part of ensuring high availability and fault tolerance in Kafka's distributed architecture.

Zookeeper-Based Leader Election (Pre-KRaft Mode)

In Kafka versions that use Zookeeper for coordination (versions before KRaft mode), the leader election process works as follows:

Zookeeper Notification: Each partition in Kafka has a leader and replicas managed through Zookeeper. When the leader broker for a partition fails or becomes unreachable, Zookeeper detects this failure by observing the absence of the leader broker in the system.

Replica Set Election: Zookeeper checks which follower replicas are available and up-to-date with the leader's log (using the log offset). Zookeeper ensures that the follower is in sync with the leader's data.

New Leader Selection: Zookeeper triggers the leader election process by selecting one of the available followers to become the new leader. The new leader must be one of the in-sync replicas (ISR) that has the most recent data (i.e., the latest log entries).

A new leader is chosen from the followers that are caught up with the leader's log (i.e., the follower that is closest to the leader in terms of the latest offset).

New Leader Assigned: Once a new leader is selected, the Kafka brokers are updated, and the new leader is now responsible for handling reads and writes for that partition.

Replication Continues: The other followers replicate data from the newly elected leader, and Kafka continues to function normally.

KRaft Mode (Kafka Raft Protocol) Leader Election (Kafka 2.8 and later)

In KRaft mode, which is a newer mode in Kafka designed to remove the dependency on Zookeeper, Kafka itself manages leader elections using the Kafka Raft Protocol. The leader election process in KRaft mode works similarly, but without involving Zookeeper.

Kafka Broker Monitoring: In KRaft mode, Kafka brokers are responsible for managing metadata and partition leadership directly. When a leader broker fails, the remaining Kafka brokers automatically detect the failure.

Leader Election: Kafka brokers initiate the leader election process by choosing the most up-to-date replica (the replica with the highest offset) from the in-sync replicas for that partition. The broker with the most recent log entries becomes the new leader.

Metadata Propagation: Kafka brokers update their internal metadata to reflect the new leader. The leader replica is assigned to handle reads and writes for the partition, while followers replicate the data.

26. How to design and banking system and explain its architecture?

A **banking system architecture** must be designed to handle secure, scalable, and reliable financial operations while ensuring regulatory compliance. Here's a high-level design covering the architecture, components, and data flow.

1. Requirements Analysis

Functional Requirements:

- **User Management:** Register, authenticate, update profiles.
- **Account Management:** Open/close accounts, check balances.
- **Transactions:** Fund transfers, withdrawals, deposits, bill payments.
- **Loan Management:** Loan applications, approvals, repayments.
- **Audit and Logging:** Maintain immutable transaction history.
- **Security:** Data encryption, secure access control, fraud detection.

Non-Functional Requirements:

- **Scalability:** Handle millions of concurrent transactions.

- **Availability:** 99.99% uptime with failover strategies.
 - **Performance:** Real-time processing with low latency.
 - **Consistency:** Strong consistency for financial transactions (ACID).
 - **Security:** GDPR compliance, PCI DSS standards.
 - **Disaster Recovery:** Multi-region backup and failover.
-

2. High-Level Architecture Overview

A modern **banking system** typically follows a **Microservices Architecture** with a layered design.

Core Layers of the Architecture:

Presentation Layer (Client Layer)

- Web App, Mobile App, ATM, Third-Party Integrations
- Manages customer interactions and authentication.

API Gateway Layer (Middleware)

- Central entry point for requests (e.g., Kong, Apigee).
- Manages routing, load balancing, rate limiting, and security.

Business Logic Layer (Service Layer)

Microservices-based design:

- **User Service:** Manages user registration and KYC.
- **Account Service:** Manages different account types.
- **Transaction Service:** Handles fund transfers and reconciliations.
- **Loan Service:** Manages loan applications and repayments.
- **Notification Service:** Sends alerts for critical events.

Data Layer (Storage and Databases)

- **Relational Databases (RDBMS):** PostgreSQL, MySQL for critical financial data.
- **NoSQL Databases:** MongoDB, Cassandra for logs and analytics.
- **Distributed Cache:** Redis for low-latency balance checks.

Event Streaming and Messaging Layer

- **Apache Kafka / RabbitMQ:** Event-driven architecture for asynchronous processing (e.g., transaction confirmations, notifications).

Security and Compliance Layer

- **OAuth2 / JWT:** Authentication and authorization.
 - **Encryption:** TLS for data in transit, AES-256 for data at rest.
 - **Fraud Detection Services:** Machine learning-based anomaly detection.
-

3. Data Flow in a Fund Transfer Example:

Scenario: User initiates a bank transfer.

Steps:

User initiates transfer from the mobile app.

API Gateway receives the request and forwards it to the **Transaction Service**.

Transaction Service performs the following:

- Validates the user and verifies balance through **Account Service**.
- Deducts the amount from the sender's account.
- Credits the amount to the receiver's account.

Event Trigger:

- Publishes an event to **Kafka** for logging and audit trails.

Notification Service:

- Sends SMS/Email confirmation to both parties.

Data Storage:

- The transaction details are stored in a **relational database**.
-

4. Key Components and Design Considerations:

Microservices Design:

- Each service handles a **single responsibility**.
 - Services can scale independently.
 - Stateless services with shared storage and messaging.
-

Database Design:

- **Relational Databases (ACID):** PostgreSQL/MySQL for core banking transactions.
 - **NoSQL Databases:** MongoDB for audit logs and historical data.
 - **Sharding and Partitioning:** For horizontal scaling of large datasets.
-

Consistency and Transactions:

- **Strong Consistency:** Required for financial transactions (ACID compliance).
 - **Two-Phase Commit (2PC):** For distributed transactions across microservices.
 - **Eventual Consistency:** For non-critical operations like notifications.
-

Security and Compliance:

- **Authentication & Authorization:** OAuth2, JWT.
 - **Data Encryption:** TLS for data in transit, AES-256 for data at rest.
 - **Compliance:** GDPR, PCI DSS, ISO 27001 standards.
-

Scalability and Fault Tolerance:

- **Load Balancing:** API Gateway with load balancing (e.g., Kong, NGINX).
 - **Replication:** Multi-region database replication for failover.
 - **Kafka:** Event-driven architecture for asynchronous event handling.
 - **Circuit Breaker:** Resilient service patterns like Hystrix.
-

5. Technology Stack:

- **Frontend:** React.js, Angular, Flutter.
 - **Backend:** Java (Spring Boot), Node.js, Python (FastAPI).
 - **Databases:** PostgreSQL, MongoDB, Cassandra.
 - **Messaging:** Kafka, RabbitMQ.
 - **Security:** OAuth2, JWT, Keycloak.
 - **Cloud:** AWS (RDS, S3, Lambda, EC2).
-

6. Summary of Services:

- **User Service:** Manages users and KYC.
- **Account Service:** Handles account management.
- **Transaction Service:** Manages transfers and balance updates.
- **Loan Service:** Loan approvals and repayments.
- **Notification Service:** Real-time alerts for transactions.

27. How to convert monolith app to micro services?

Converting a **monolithic** application to **microservices** is a complex but rewarding process that allows for greater scalability, flexibility, and maintainability. The key challenge is breaking down the tightly-coupled monolithic structure into smaller, loosely-coupled services that are independently deployable and manageable.

Here are the **steps to convert a monolithic application to microservices**:

1. Evaluate the Monolith

- **Understand the existing system:** Start by thoroughly understanding the monolithic application. Identify the key components and how they interact with each other (databases, APIs, internal logic, etc.).
- **Identify business domains:** Group the application functionality into logical business domains or bounded contexts (e.g., user management, billing, inventory). This will help in defining the boundaries of your future microservices.

2. Define Microservices Boundaries

- **Break down the monolith:** Identify independent modules or components that can be isolated as microservices. Each microservice should represent a specific business function or domain.
- **Apply Domain-Driven Design (DDD):** Use principles from DDD to define bounded contexts, which can help to identify services that have clear responsibilities and interfaces.
- **Decide on the granularity:** Decide whether the services will be fine-grained (i.e., smaller services) or coarse-grained (i.e., fewer larger services), based on the complexity of the business logic.

3. Choose the Right Technology Stack

- **Select appropriate technologies:** Decide whether to keep the existing technology stack or adopt new technologies for each microservice (e.g., different languages, frameworks, databases).
- **Consider APIs:** Each microservice should expose well-defined APIs, typically using RESTful services or gRPC for communication.
- **Data storage:** Microservices should ideally have their own databases, so you may need to migrate from a shared database to multiple isolated data stores. Consider whether to use SQL, NoSQL, or hybrid storage solutions.

4. Start with One Service (Strangler Pattern)

- **Strangler pattern:** Rather than rewriting the entire monolith at once, start by gradually replacing parts of the monolith with microservices. You can implement new features or re-implement existing features as microservices while keeping the monolithic application running.

- **Isolate and migrate incrementally:** Migrate each functional module from the monolith to its own microservice one at a time. As you migrate, ensure the new microservice does not depend on the monolithic codebase.

5. Define Inter-Service Communication

- **Communication patterns:** Microservices communicate with each other via lightweight protocols, typically over HTTP/REST, gRPC, or messaging queues (e.g., Kafka, RabbitMQ). Ensure there's an efficient and reliable communication pattern.
- **Synchronous vs. Asynchronous:** Determine whether the communication will be synchronous (e.g., REST API calls) or asynchronous (e.g., event-driven with message queues).
- **API Gateway:** Implement an API Gateway to aggregate multiple microservices into a single entry point for clients and handle routing, security, load balancing, and more.

6. Implement a Service Registry and Discovery

- **Service Registry:** Use a service registry (e.g., Consul, Eureka, or Kubernetes) to keep track of the instances of each microservice, which can be dynamically discovered and accessed by other services.
- **Dynamic Scaling:** With microservices, services can scale independently. Use load balancing and service discovery to dynamically scale microservices based on demand.

7. Manage Data Consistency

- **Event-driven architecture:** Microservices often require eventual consistency rather than strict ACID transactions. Use event-driven architectures to synchronize data across services (e.g., publish events using Kafka or RabbitMQ).
- **Database per service:** Each microservice should have its own database to prevent tight coupling. This requires strategies like event sourcing, CQRS (Command Query Responsibility Segregation), and compensating transactions for eventual consistency.

8. Implement CI/CD (Continuous Integration/Continuous Deployment)

- **CI/CD pipelines:** Automate testing, building, and deployment for each microservice. Tools like Jenkins, GitLab CI, and CircleCI can help implement CI/CD pipelines.
- **Containerization:** Use Docker to containerize each microservice. This makes it easier to deploy, scale, and maintain them independently.
- **Orchestration:** Use orchestration platforms like **Kubernetes** or **Docker Swarm** to manage, scale, and deploy microservices in a distributed manner.

9. Monitor and Manage Microservices

- **Logging and Monitoring:** Implement centralized logging and monitoring tools like **ELK stack** (Elasticsearch, Logstash, Kibana), **Prometheus**, and **Grafana**. This helps in tracking each service's performance and debugging issues.
- **Distributed Tracing:** Use tools like **Jaeger** or **Zipkin** for tracing requests across services, which is especially useful for debugging and identifying performance bottlenecks.
- **Error Handling and Resilience:** Implement circuit breakers (e.g., using **Hystrix** or **Resilience4j**) to gracefully handle service failures and retries.

10. Security and Governance

- **Service Security:** Ensure secure communication between services using **TLS** encryption, OAuth, or API keys.
- **Authentication and Authorization:** Use a centralized authentication system (e.g., OAuth, OpenID Connect) to handle user authentication and authorization across microservices.
- **API Management:** Implement an **API Gateway** to centralize access control, rate limiting, and security checks.

11. Test and Deploy Microservices

- **Testing:** Ensure robust testing for each microservice. Unit testing, integration testing, and end-to-end testing are critical for microservices.
- **Deployment Strategy:** Use blue/green deployments or canary releases to deploy microservices with minimal downtime. Containers and orchestration platforms (e.g., Kubernetes) help in managing deployments at scale.

12. Iterate and Optimize

- **Refactor as needed:** After breaking down the monolith and creating microservices, continually monitor and optimize both the architecture and individual services.
 - **Manage technical debt:** Keep an eye on the growing complexity of maintaining multiple microservices. Introduce necessary refactoring to ensure the system remains maintainable in the long run.
-

Challenges of Converting a Monolith to Microservices

- **Increased Complexity:** Microservices introduce complexity in terms of communication, data management, deployment, and monitoring.
 - **Data Consistency:** Managing data consistency across distributed services can be tricky, and special patterns like event sourcing or eventual consistency need to be used.
 - **Deployment and Scaling:** While microservices are designed for scalability, managing deployment and scaling of many services can be difficult without proper tools and infrastructure (e.g., Kubernetes).
 - **Distributed Debugging:** Debugging issues across many independent services can be more challenging than in a monolithic application.
-

Best Place to learn System Design: [gopi582/awesome-system-design-resources: Learn System Design concepts and prepare for interviews using free resources.](#)

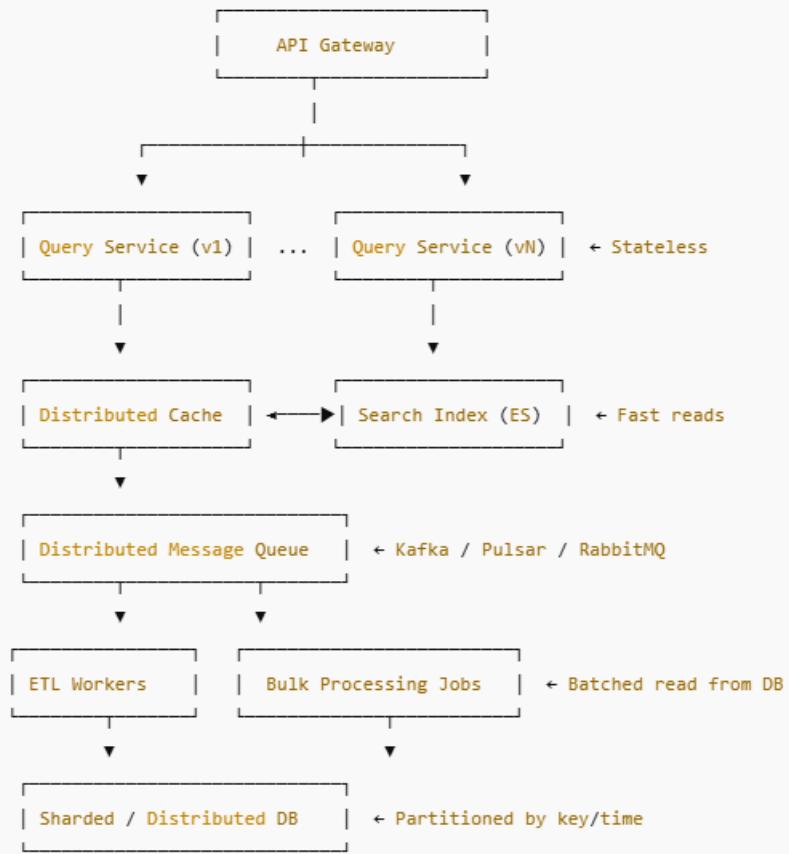
System design interview questions: [System Design Interview Questions 2025 | System Design Interview Questions & Answers | Intellipaat](#)

28. Design micro services architecture to pull billions of records from database with out latency and availability?

✓ High-Level Architecture Overview

mathematica

Copy Edit



Key Design Components

1. Sharded/Partitioned Database

- Use **horizontal partitioning** (e.g., by customer ID, time, region).
- Use databases built for **massive scale**, e.g.:
 - ScyllaDB, Cassandra, ClickHouse, Google BigQuery, or Amazon Redshift.
 - Or use PostgreSQL with **Citus** extension for distribution.

2. Data Ingestion Layer

- Use **Kafka** or **Apache Pulsar** to decouple read/write loads.
- Break large queries into **smaller tasks** that can be **distributed** to workers.

3. ETL/Worker Services

- Asynchronous **batch readers** that pull data from the database.
- Use **Spring Boot + Kafka consumers** with **batch processing**.
- Workers are stateless and **scalable horizontally**.

4. Distributed Cache Layer

- Use **Redis**, Memcached, or **Aerospike** for fast access.
- Cache **hot queries**, **pagination cursors**, and **frequent filters**.

5. Search Indexing

- Index structured data in **Elasticsearch** for fast, paginated search.
- Use **data sync pipelines** (Kafka → Elasticsearch).

6. Query Services (Stateless Microservices)

- Expose **REST/GraphQL APIs** for clients.
- Read from cache, search index, or trigger background fetch via message queue.
- Use **WebFlux** (non-blocking) or **gRPC** for fast async APIs.

7. API Gateway

- Rate limiting, load balancing, auth.
- Can use **Istio**, **Kong**, or **Spring Cloud Gateway**.

8. High Availability

- Use **Kubernetes** for autoscaling and container orchestration.
- Multi-region deployment with **geo-redundant databases**.
- Use **circuit breakers** (e.g., Resilience4j) to isolate failures.

9. Monitoring and Observability

- Metrics: **Prometheus + Grafana**.
- Logs: **ELK stack** or **Loki**.
- Tracing: **Jaeger** or **OpenTelemetry**.

Performance Optimization Tips

| Concern | Solution |
|------------------------|--|
| High Read Latency | Cache responses, use async non-blocking I/O (e.g., WebClient, Netty) |
| Billions of Records | Partitioned reads, distributed query execution |
| Data Freshness | Kafka CDC pipelines, Redis TTL, async refresh |
| Rate Limiting | API Gateway with user-level throttling |
| Fault Tolerance | Retry + Circuit Breaker + Bulkhead pattern |
| Pagination Performance | Use cursor-based pagination , not OFFSET |

Example Tech Stack

| Layer | Tools/Tech |
|---------------|---------------------------------|
| API Layer | Spring Boot, gRPC, GraphQL |
| Gateway | Kong, Spring Cloud Gateway |
| Message Queue | Apache Kafka, Pulsar |
| Data Store | Cassandra, ScyllaDB, ClickHouse |
| Search Index | Elasticsearch |
| Cache | Redis, Aerospike |
| Orchestration | Kubernetes + Helm |
| Monitoring | Prometheus, Grafana, Jaeger |

Example Use Case: Pull Last 10 Billion Logs

1. User calls `/logs?start=2024-01-01&end=2025-01-01`.
2. API checks if query is cached in Redis.
3. If not:
 - Sends request to Kafka: `"fetch logs between time range"`.
 - Background ETL service processes partitioned query.
 - Streams results into Elasticsearch.
4. Once ready, user is notified (polling or WebSocket).
5. Query service streams from Elasticsearch to user.